Milestone 2 Report

The intended audience/users for this project are people who would like to pick up on reading, long-time readers, consumers, editors who are interested in the topics of books, publishers, booksellers, and authors.

For task 1, with the use of a normalized stacked bar chart, the audience can identify successful publication years by comparing the height of "more than 3 awards" categories bar over the years and booksellers can develop marketing strategies to promote books based on most successful years and so on.

For task 2, with the use of a line chart, the audience can quickly identify which genre is most popular over the years and learn which types of books are in high demand and potentially identify gaps in the market.

For task 3, with the use of heatmap, users can base on colors to quickly identify which publishers are dominant in what genres and identify potential gaps in the market.

For task 4, the audience can find out which genre of books are popular across Amazon and Goodreads and if books' ability to stay on the best-selling list is correlated with popularity.

For task 5, viewers can find books they are looking for by dynamic queries (year, genre) and make informed decisions about where to buy books to potentially save money on their purchase.

For task 6, this view can help users identify which books are highly rated but have a low number of reviews, potentially find a hidden gem that could become popular with readers or identify books that are popular (have lots of reviews) due to the good marketing but actually are not that good (low ratings).

Task 7, this view helps users to understand the relationship between critical acclaim and commercial success by looking at the distribution of number of weeks on best-selling list wrt number of awards in box plot. Users can also somewhat get the big picture of correlation between the two attributes if it exists.

Tasks:

Task 1: Which publication year has the highest proportions of books that have more than 3 awards?

Task 2: What are the Most Popular Book Genres Over the Years?

Task 3: For the Top 10 most popular genres, how many books were published by Top 5 publishers?

Task 4: Do fiction books stay on the NYT best-selling list longer than non-fiction books and are they more popular on Amazon and Goodreads?
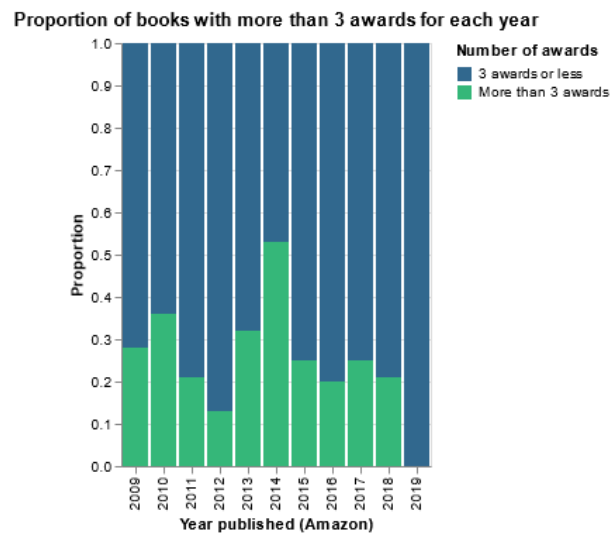
Task 5: How do the prices differ between Goodreads and Amazon platforms per book listed in Fiction/Non-fiction genre?

Task 6: What is the number of books per rating filtered by number of reviews on Amazon and Goodreads?
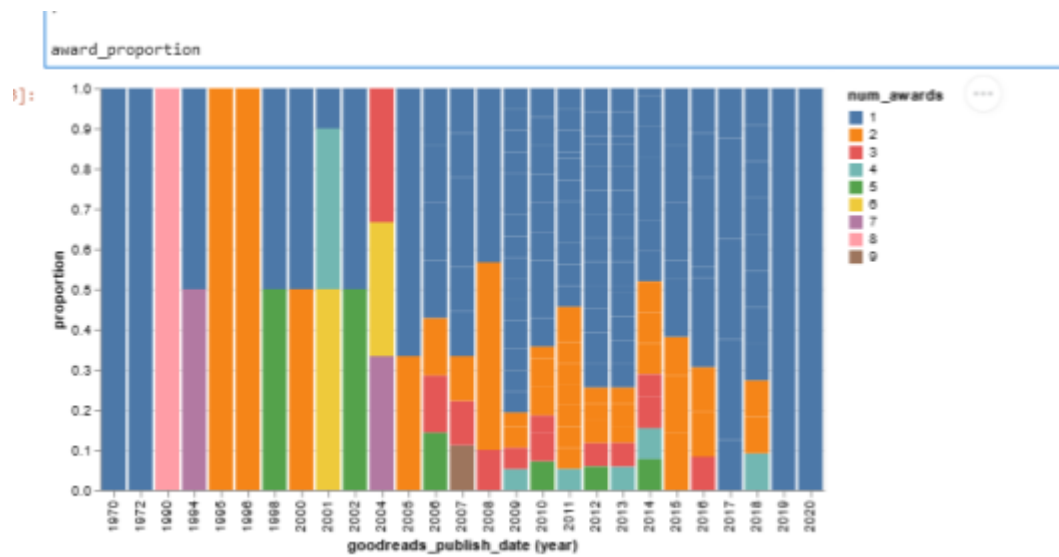
Task 7: Do books with more awards tend to stay longer on the NYT best selling list?

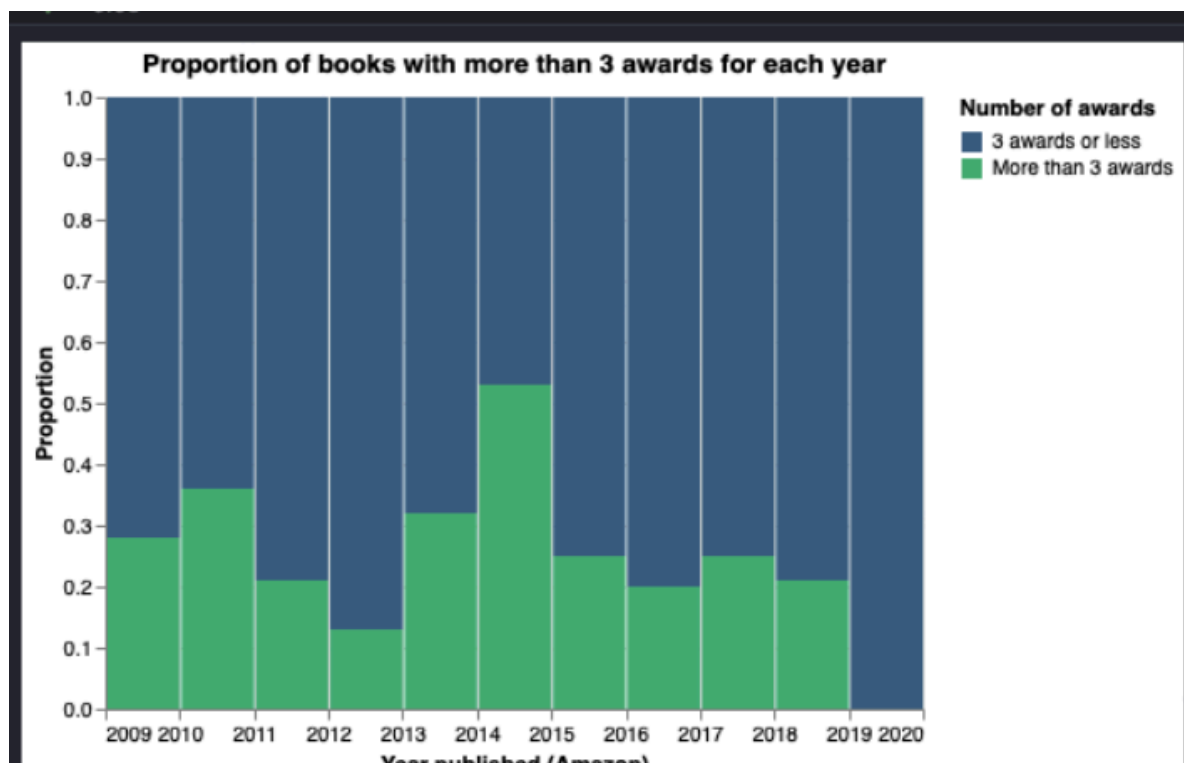View 1: ´Proportions of books with more than 3 awards for each year
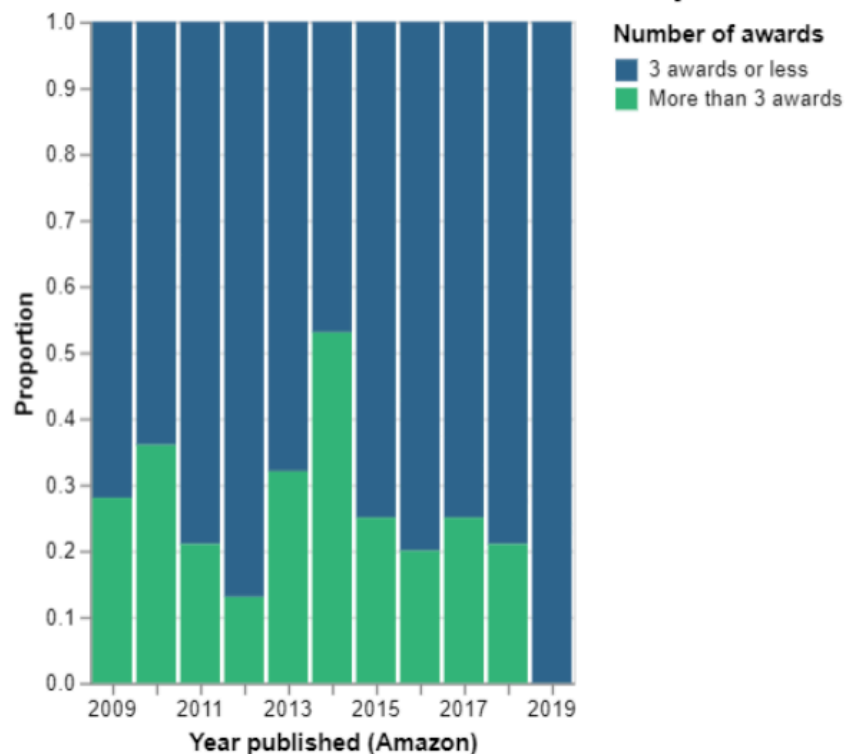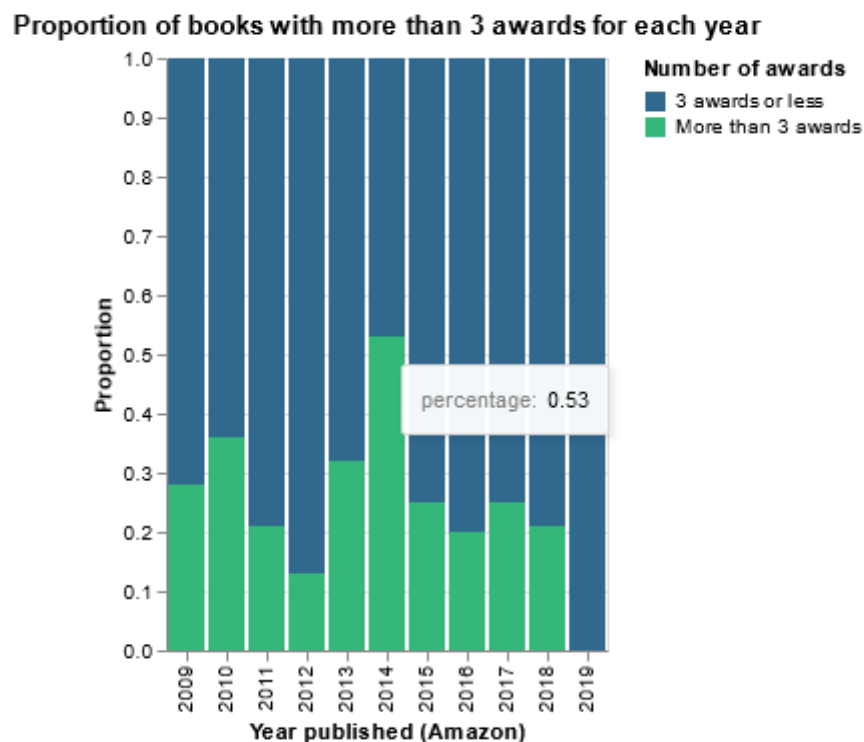
- Image of visualization:



- Include screenshots of previous iterations of the visualization

**Proportion of books with more than 3 awards for each year**



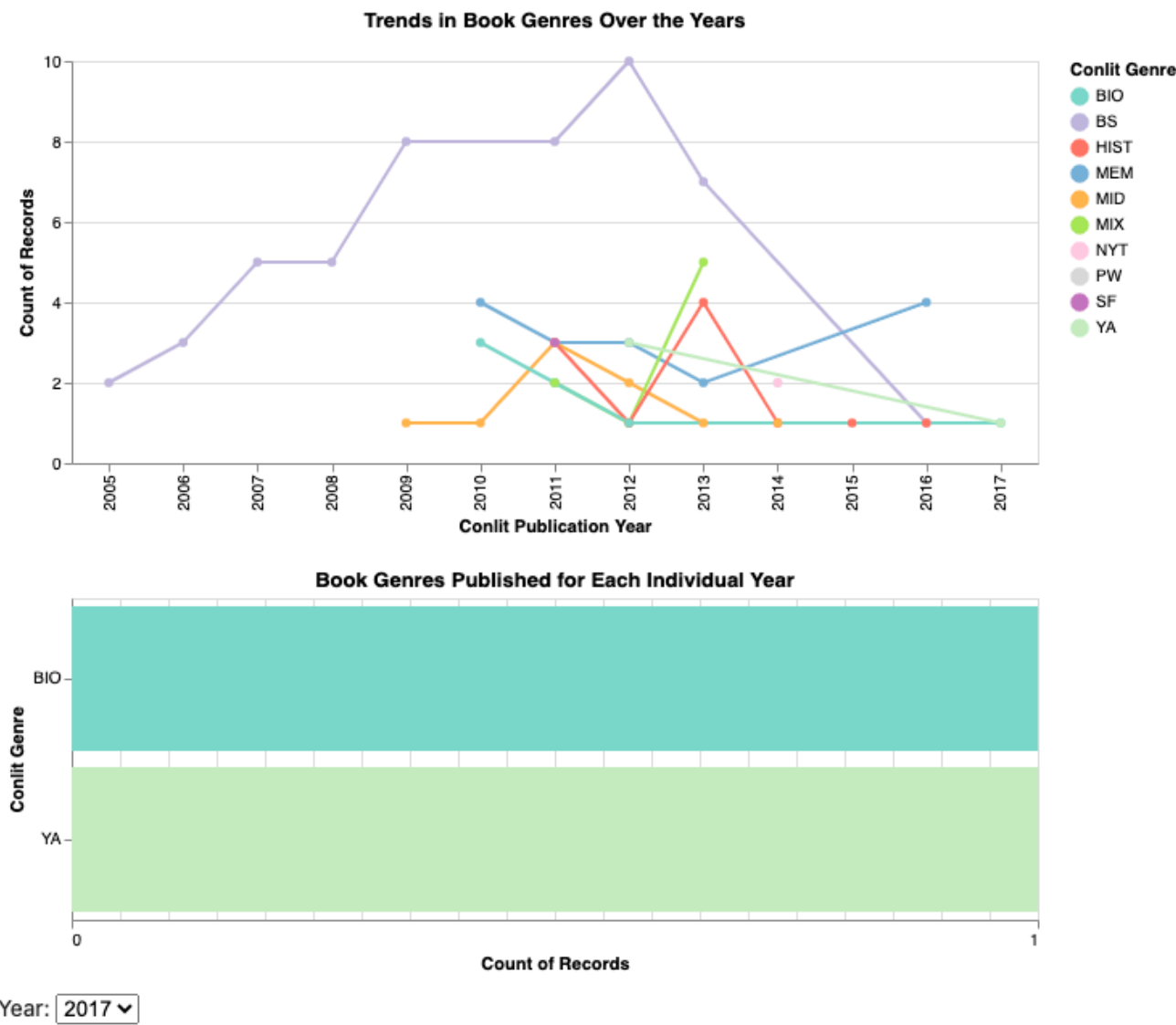**Proportion of books with more than 3 awards for each year**

- Task addressed:
  - Which year has the highest proportions of books that have more than 3 awards?
- Explanation:
  - Marks: The stacked bars that represent the proportion of books with more than 3 awards for each year.
  - Channels:
    - The x-axis: The x-axis represents the year when it was published on Amazon
    - The y-axis: The height of each bar chart is determined by the percentage field.
    - Color hue: The color is determined by "num_awards" encoded into the color channel
    - Tooltip: Shows the percentage when mouse is hovered over a bar
- How channels were exploited:
  - Color hue chosen over color saturation or luminance because it is a categorical variable
  - Color channel has high discriminability in this view as we can easily see two levels of attribute for the number of awards. Moreover, it is also friendly for color-blind users.
- Describe the interaction
  - The tooltip states the exact percentage for each stack in the stacked chart
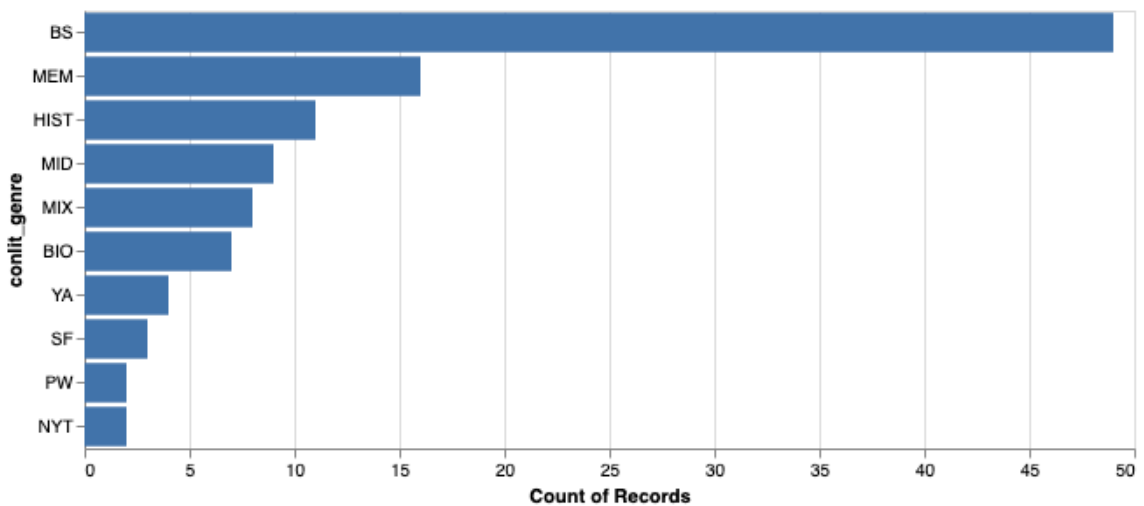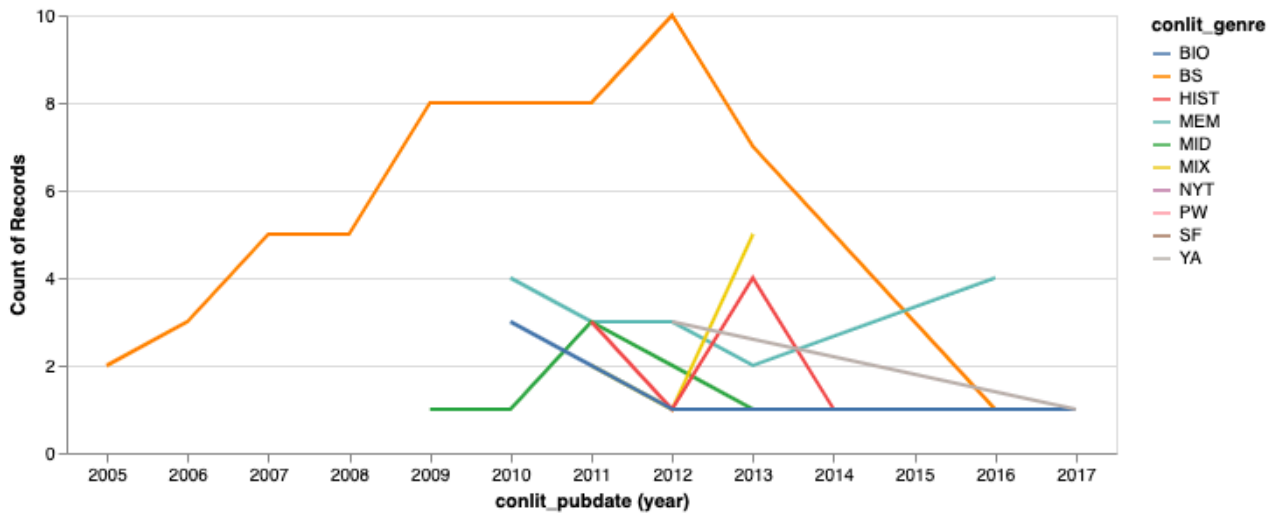


Proportion of books with more than 3 awards for each year

- ○ Characteristics of Interaction and interactivity
  - ■ Hover:
  - ■ Action: select
  - ■ Event type: hover
  - ■ Reaction type: highlight
  - ■ Interaction coupling: unidirectional
  - ■ Action Elements - Focus: direct, Presence: implicit, Granularity: atomic
  - ■ Reaction Elements - Activation: Immediate, Spread: Self contained form, Flow: Discrete
- ○ Critique the view
  - ■ The choice of bar chart and line chart is very straightforward and easy to understand for a wide range of people, which is suitable for our target audience.
  - ■ Bar chart uses a common axis which increases accuracy for comparison as users only need to compare the heights of the bars from one another.
  - ■ Since we used proportion we cannot access the count directly without any interaction and we can't tell how many of those books won exactly 4 awards or 5 awards and so on.
  - ■ There is no visual cue that highlights the lowest or highest proportions of highly awarded books.
  - ■ The tooltip could be improved upon by including both of the percentages of the books with "3 ≥ awards"  and "more than 3 awards" to allow a precise comparison when hovering over the bars.
  - ■ Color channel has high discriminability in this view as we can easily see two levels of attribute for the number of awards. Moreover, it is also friendly for color-blind users.
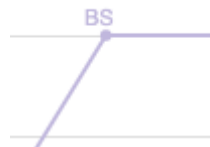
View 2:

- Image of visualization



**Trends in Book Genres Over the Years**

**Book Genres Published for Each Individual Year**

Year: 2017 ▾

- Include screenshots of previous iterations of the visualization

- Task addressed: What are the most popular book genres over the years?
- Explanation:
  - Marks:
    - Line Chart: Points and Lines represent each genre in each year
    - Bar Chart: Lines representing the genres for a specific year specified by the drop-down
  - Channels
    - X and Y Position on an aligned scale encodes the count of each genre
    - Color hue encodes the genre
  - Characteristics of Channels that were exploited
    - Color hue is used because genre is a categorical variable
    - X and Y position, so we have high accuracy since we want to compare the counts

- We made sure to choose colors that are color blind friendly so it's accessible to a wide range of audiences
- Colors in the bar graph matches colors in the line graph so it's easier to map the genres to the other visualization
- Describe the interaction
  - Hover over the points to show the genre



  - Click on the legend so the line chart highlights the line with the chosen genre



  - Dropdown specifies which year to show in the bar chart



- Characteristics of Interaction and interactivity

  Hover:

  - Action: select
  - Event type: hover
  - Reaction type: highlight
  - Interaction coupling: unidirectional
  - Action Elements - Focus: direct, Presence: implicit, Granularity: atomic

- - - ■ Reaction Elements - Activation: Immediate, Spread: Self contained form, Flow: Discrete
  - Click:

    - ■ Action: select
    - ■ Event type: click
    - ■ Reaction type: highlight
    - ■ Interaction coupling: unidirectional
    - ■ Action elements - Focus: indirect, Presence: implicit, Granularity: atomic
    - ■ Reaction Elements - Activation: immediate, Spread: self contained form, flow: discrete

  - Dropdown:

    - ■ Action: select
    - ■ Event type: click
    - ■ View: Partition
    - ■ Views Data: share data
    - ■ reaction type: change data
    - ■ Interaction Coupling: unidirectional
    - ■ Action Elements - Focus: indirect, Presence: Explicit, Granularity: Composite
    - ■ Reaction Elements - Activation: Immediate, Spread: self-contained form, Flow: discrete
    - ■ Affordances: upside down triangle indicate a drop-down menu with choices available to pick from
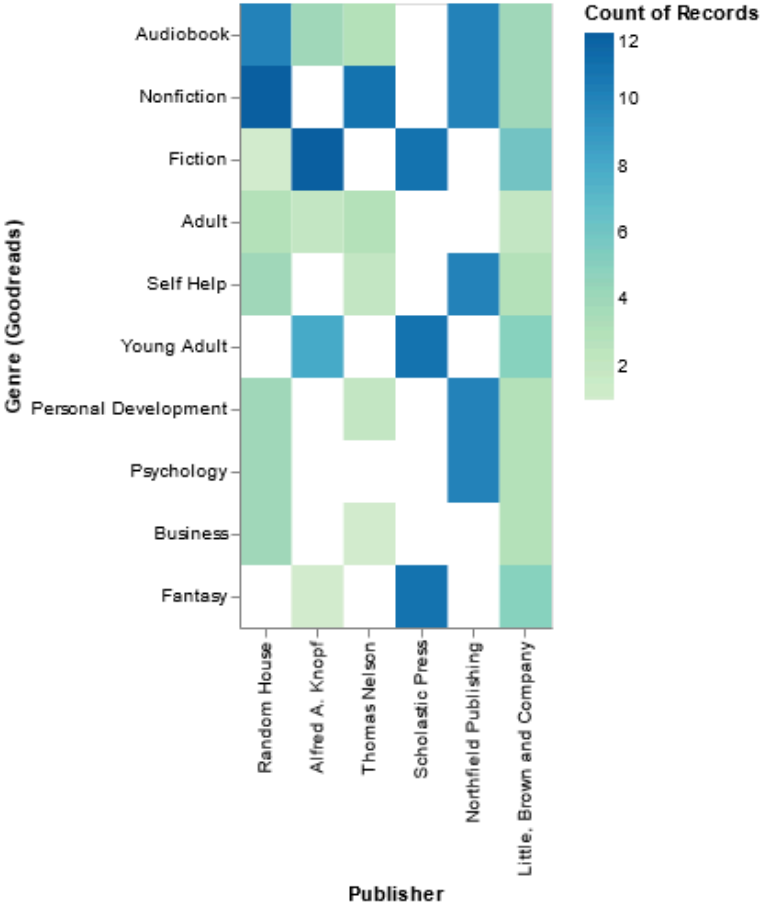  - ○ Critique
    - ■ The choice of bar chart and line chart is very straightforward and easy to understand for a wide range of people, which is suitable for our target audience.
    - ■ Amount of colors used is not ideal but since the lines overlap with one another it may be more confusing to only use one color for all the lines. The colors are also different enough that there is sufficient discriminability between the genres.
    - ■ The color palette is catered to work for the most common form of color blindness, deuteranomaly. It may not necessarily work for other forms of color blindness.
    - ■ Using color and hover to encode genre may be redundant but since the lines overlap with each other it helps identify the points easier.

- Using color in the bar graph can be redundant since the genre is already stated in the axis but matching the color in the bar graph to the color in the line chart helps connect the two visualizations
- Bar chart uses a common axis which increases accuracy for comparison as users only need to compare the heights of the bars.
- Abbreviated genres requires the audience to search for the meaning of the abbreviation if it isn't obvious to them
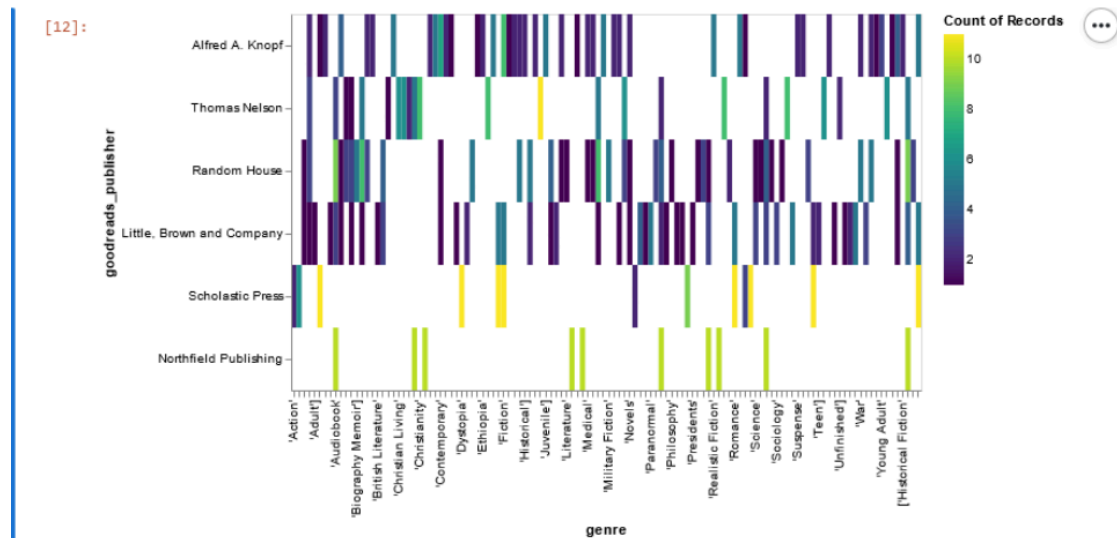
View 3:

- Image of visualization



For Top 10 most popular genres, how many books were published by Top 5 publishers

- Include screenshots of previous iterations of the visualization

heatmap





- Task addressed:
  - For the top 10 most popular genres, how many books were published by top 5 publishers?
- Explanation:
  - **Marks**: Rectangular bars
  - Channels:
    - X-axis: represents the publishers sorted by popularity.
    - Y-axis: represents the genres sorted by popularity.
    - Color saturation/luminance: represents the count of books published for each combination of genre and publisher. The color scheme is

sequential from a range from light green to dark blue, the saturation and luminance change as the value of counts increases and the intensity of color also does
- ■ Tooltip: Provides additional information and shows the specific count of books published for each combination of genre and publisher when hovering over the corresponding bar.
- ○ Characteristics of Channels that were exploited
  - ■ X and Y channels both use categorical data to represent the publishers and genres, respectively.
  - ■ Color channel uses numerical data to represent the count of books published by the top 6 publishers of each genre, the color gradient varies in color intensity determined by the count value.
- ○ Describe the interaction
  - ■ When the mouse is hovered over each rectangle in the heatmap, the tooltip outputs the count of published books for each combination of genre and publisher.



]: **For Top 10 most popular genres, how many books were published by Top 5 publishers**

- ○ Characteristics of Interaction and interactivity
  - ■ Hover
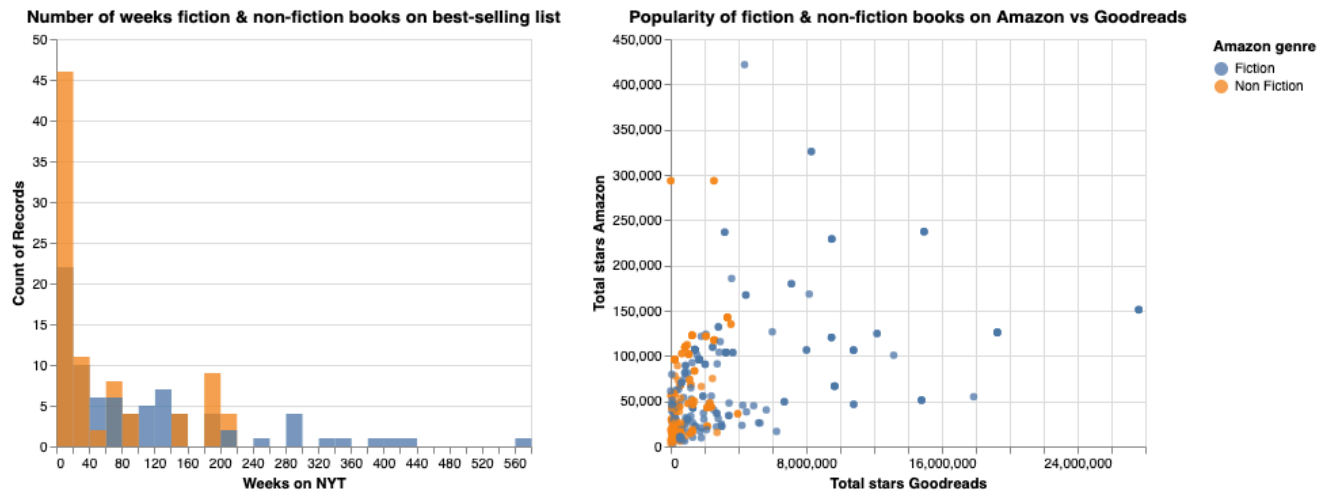    - ● Action: Select
    - ● Event type: hover

- Reaction type: highlight
- Interaction coupling: unidirectional
- Action Elements - Focus: direct, Presence: implicit, Granularity: atomic
- Reaction Elements - Activation: Immediate, Spread: Self contained form, Flow: Discrete
  - Critique
    - The empty rectangles in the visualization may distract from the insights trying to be communicated
    - The data is not very sufficient, I think if more data were to be collected, the view could deliver a much more meaningful insight.
    - Color channel's discriminability is good enough for our task as we would want to look into (genre, publisher) with the most count - darkest color, which is very distinctive from all other colors. Moreover, I believe this color choice is friendly for users who are colorblind as in our case, the darkest color would be instinctively eye-catching expressing data with most counts.

View 4:

Task: Do fiction books stay on the NYT best-selling list longer than non-fiction books, and are they more popular on Amazon and Goodreads?



- Image of visualization
- Include screenshots of previous iterations of the visualization

  Previous iterations of the visualization have the same charts as above, the difference is the interactions implemented.

- Task addressed:

  The bar chart addresses whether fiction books stay on the NYT best-selling list longer than non-fiction books.

  Scatter plot addresses whether fiction books are more popular than non-fiction books on Amazon and Goodreads.

  Using interaction interval selection on bar chart to choose books on scatter plot based on the number of weeks stay on the best-selling list addresses the underlying question whether fiction books that stay on the best-selling list for longer times are more popular on Amazon and Goodreads.

- Explanation:
  - Marks

    Bar chart - bar: count of books per number of weeks on NYT best-selling list.

    Scatter plot - circle: books.

  - Channels

Bar chart: Number of weeks on NYT best-selling list is encoded on x-axis, count of books is encoded on y-axis. Moreover, Amazon genre (fiction/non-fiction) is encoded as color hue.

Scatter plot: Total stars on Goodreads is encoded on x-axis, total stars on Amazon is encoded on y-axis. Moreover, Amazon genre (fiction/non-fiction) is encoded as color hue.

○ Characteristics of Channels that were exploited
  - Bar chart:

For x-axis, the scale runs from 0 to 580 (weeks).

For y-axis, the scale runs from 0 to 50 (count of books).

Since bar chart uses a common axis, it increases the accuracy of the visualization in terms of comparison between genres with respect to count of books on weeks on NYT best-selling list. Users only need to compare the height of blue bar and orange bar at same x-coordinate to know which genre has more books for said number of weeks.

With the use of color hue in both charts, we can easily differentiate which genre that datapoints/books belong to.

Color hue is the best choice in this view for encoding 2 genres/categories.

For color channel, it is encoded as color hue with two levels: blue for fiction genre and orange for non-fiction genre.

  - Scatter plot: total stars = number of reviews * rating

For the x-axis, the scale runs from 0 to 28000000 (total stars).

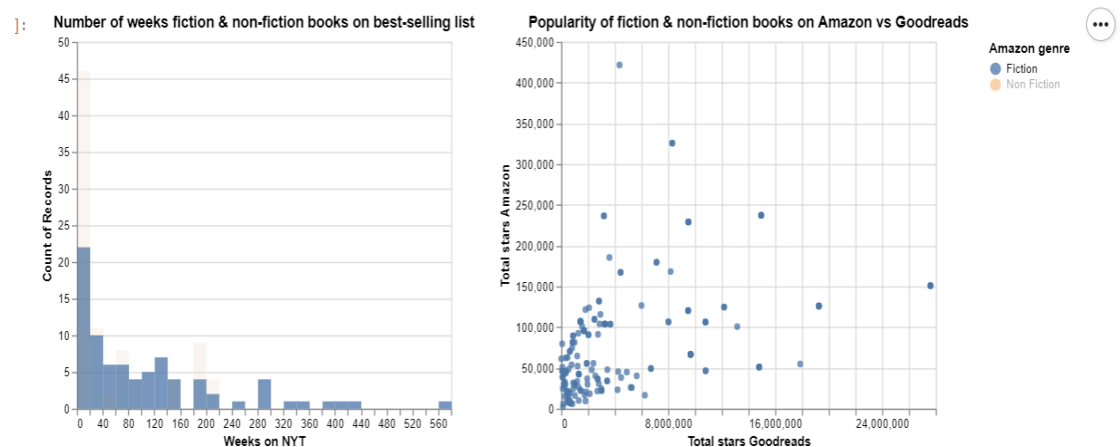For the y-axis, the scale runs from 0 to 450000 (total stars).

For color channel, it is encoded as color hue with two levels: blue for fiction genre and orange for non-fiction genre.

Scatter plot uses common axis, which is easy for users to navigate which books' genre are popular on Goodreads/Amazon by looking at datapoints' color on end of the spectrum of x-axis and y-axis.
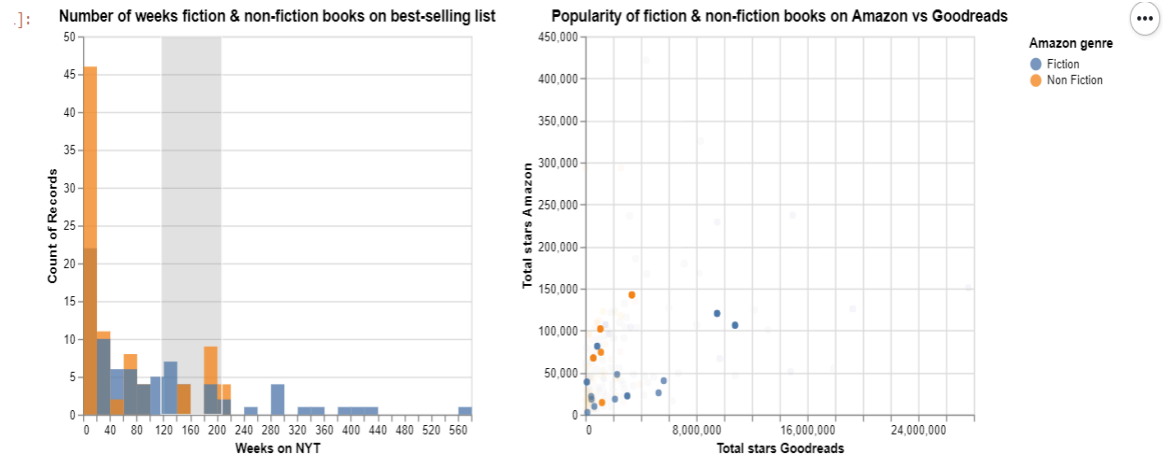
○ Describe the interaction

1. Selection (type: on click, semantics: single and multiple): on interactive legend

Selection is encoded on color legend (encoded for fiction/non-fiction genre) so that if users click on either or both genres, the corresponding data points' visual encodings will change accordingly to highlight the chosen genres' data points in bar chart and scatter plot. To be more specific, in the bar chart, the opacity encoding of bars will increase (from 0.6 to 0.75) if chosen and decrease (from 0.6 to 0.05) if not chosen. In the scatter plot, the size of circles will stay the same (30) if chosen and decrease to 0 if not chosen.



2. Selection (constrained on y-channel of bar chart, type: on shift-click, on mouse scroll, on click, semantics: null, group)

Users use this rectangular selection on bar chart to change the visual encoding of data points in the scatter plot. Opacity encoding of data points in scatter plot is increased to 0.75 if chosen in selection in bar chart and decreased to 0.01 otherwise. This allows the user to see the popularity of books in each genre on Amazon and Goodreads given how long they stay on the NYT best-selling list.

.]:



3. Selection & filter (constrained on x-channel and y-channel of scatter plot, type: on shift-click, on mouse scroll, on click, semantics: null, group)

Users use interval selection on scatter plot to filter the contents of a bar chart. This allows the user to see how long books in each genre stay on the NYT best-selling list within the selection of popularity on Amazon and Goodreads.



- ○ Characteristics of Interaction and interactivity

    1. Action: select

    Event type: click, multiple click

    Reaction type: highlight (visual encoding)

    Views: juxtapose

Spread: propagated form - bidirectional. (affected both visualizations)

Views data: share encoding, share data.

Interactivity: Focus: direct focus. Presence: implicit. Granularity: atomic

Activation: immediate

Flow: discrete flow.

2. Action: select

Event type: on shift-click, on mouse scroll, on click

Reaction type: highlight (visual encoding)

Views: juxtapose

Spread: self-contained form - unidirectional. (affect visual encoding of scatter plot only, bar chart is unaffected)

Views data: share data, share encoding.

Interactivity: Focus: direct focus. Presence: implicit. Granularity: atomic

Activation: immediate

Flow: discrete flow.

3. Action: select & filter

Event type: on shift-click, on mouse scroll, on click

Reaction type: filter data

Views: juxtapose

Spread: self-contained form - unidirectional. (affect data of bar chart only, scatter plot is unaffected)

Views data: share encoding, share data.

Interactivity: Focus: direct focus. Presence: implicit. Granularity: atomic

Activation: immediate
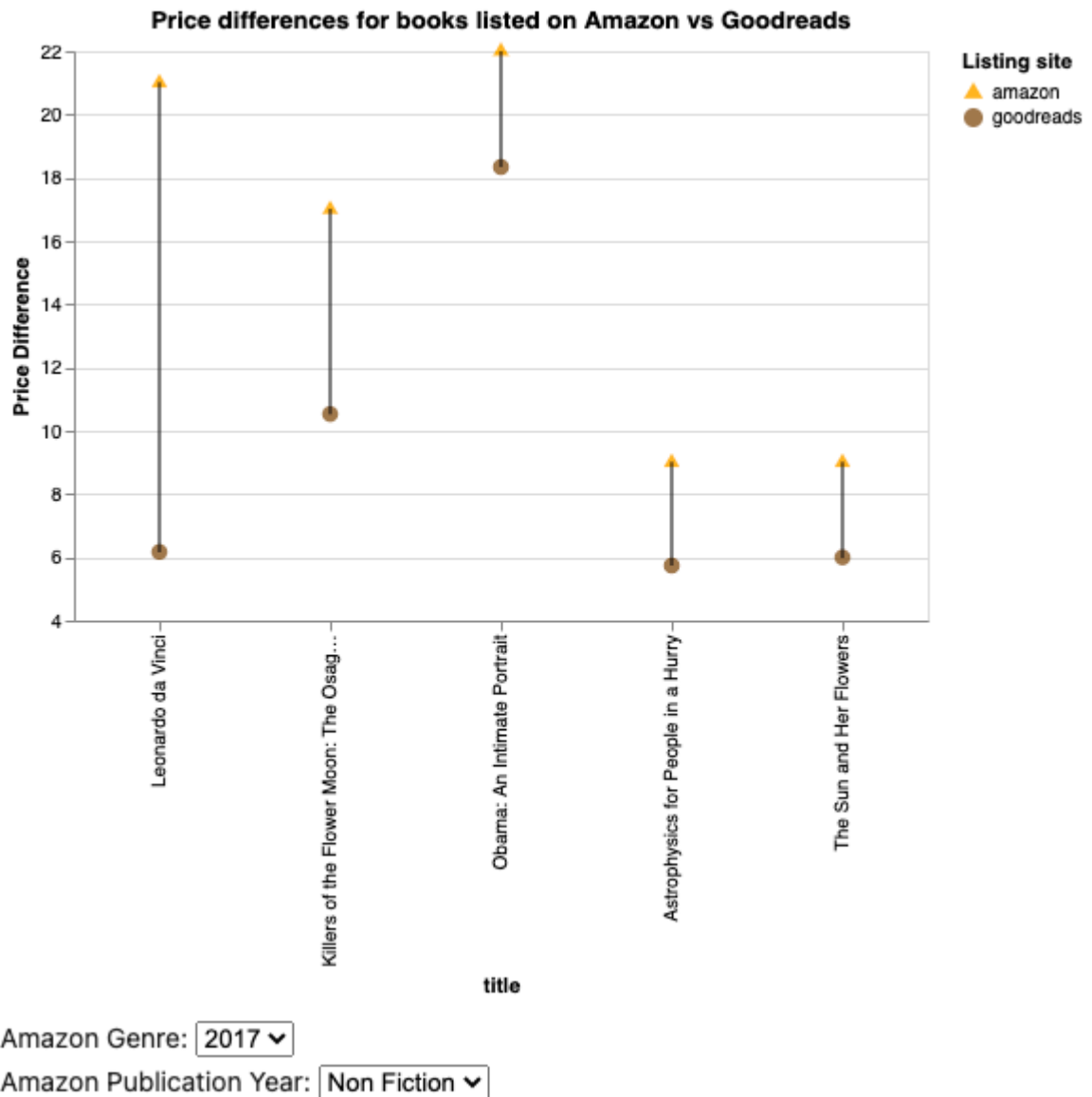
Flow: discrete flow.

- ○ Critique

With the use of color hue in both charts, we can easily differentiate which genre that datapoints/books belong to. Since the bar chart uses a common axis, it increases the accuracy of the visualization in terms of comparison between genres with respect to count of books on weeks on NYT best-selling list. Users only need to compare the height of blue bar and orange bar at the same x-coordinate to know which genre has more books for said number of weeks. This answers whether fiction books stay on the best-selling list longer than non-fiction books. Moreover, the scatter plot also uses common axis, so it is easy to navigate which books' genre are popular on Goodreads by looking at datapoints' color to the right-end of x-axis, which books' genre are popular on Amazon by looking at datapoints' color to the upper-end of y-axis, and which books' genre are popular on both platforms by looking at data points' color on upper right side of the plot. Color channel has high discriminability in this view as we can easily see two levels of attribute amazon genre and the color choice is friendly for color blind users. Lastly, the use of channel color and position is fully separable in both charts.

Furthermore, since quite a few bars overlap and circles are concentrated mostly in one place in scatter plot, the use of genre selection on interactive legend makes it easier for users to get accurate count of books per each genre in bar chart and to investigate the popularity of each genre on Goodreads and Amazon.  Additionally, the use of rectangular selection on bar chart to highlight corresponding books popularity on scatter plot answers the underlying question whether fiction books that stay on best-selling list longer are more popular as users can select bars that are on the right end of the x-axis to select fictions books that stay on the list for longer times and investigate its popularity over in the scatter plot. The use of interval selection on scatter plot provides the user with freedom to engage in further analysis of data (ie: for popular books on Amazon/Goodreads in each genre, how long do they stay on best-selling list).

To sum up, the view efficiently answers the task and the choice of charts are easy to understand for a wide range of audiences. Moreover, with the use of interactions in the view, users are provided with more freedom to adjust the amount of data they want to investigate and makes them more engaged in the data analysis process.

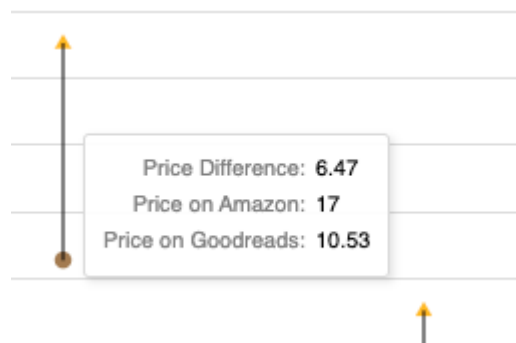View 5: Price differences for books listed on Amazon vs. Goodreads
- Image of visualization



- Include screenshots of previous iterations of the visualization
- Task addressed:
  - How do the prices differ between Goodreads and Amazon platforms per book listed in Fiction/Non-fiction genre(s)?
- Explanation:
  - Marks:
    - Points: Prices in each platform
    - Lines: Price difference
  - Channels and justification

- X-axis: Encodes book titles sorted in alphabetical order.
- Y-axis: Point marks on linear scale with non-zero starting point
  - Scatterplot: Encodes book prices for the scatter plot with the column "value" represents the price
  - Line graph: Goodreads price and Amazon price
- Shape: Triangle and circle for Amazon and Goodreads respectively.
- Color: Different color hues to also distinguish the platforms, the encoded colors represent the companies' colors used in their logos.
- Tooltip: Displays price difference between platforms when user hovers over the points.

- Exploited channel characteristics:
  - X-axis: Book titles in alphabetical order
  - Y-axis: Non-zero starting point linear scale allows easy comparison of prices between the two platforms.
  - Shape: Different shapes to distinguish each platform, triangle and circle shapes allow easy recognition.

- Describe the interaction
  - Hovering over the lines shows the price difference, and the price for each site.



Price Difference: **6.47**
Price on Amazon: **17**
Price on Goodreads: **10.53**

  - A drop-down list for both 'amazon_year', and 'amazon_genre', amazon_year represents the year it was published on amazon, and it goes through a range from 2009-2019, and the selected genres to compare in 'amazon_genre" are Non-Fiction and Fiction books

Amazon Genre: [ 2017 ✔ ]
Amazon Publication Year: [ Non Fiction ✔ ]
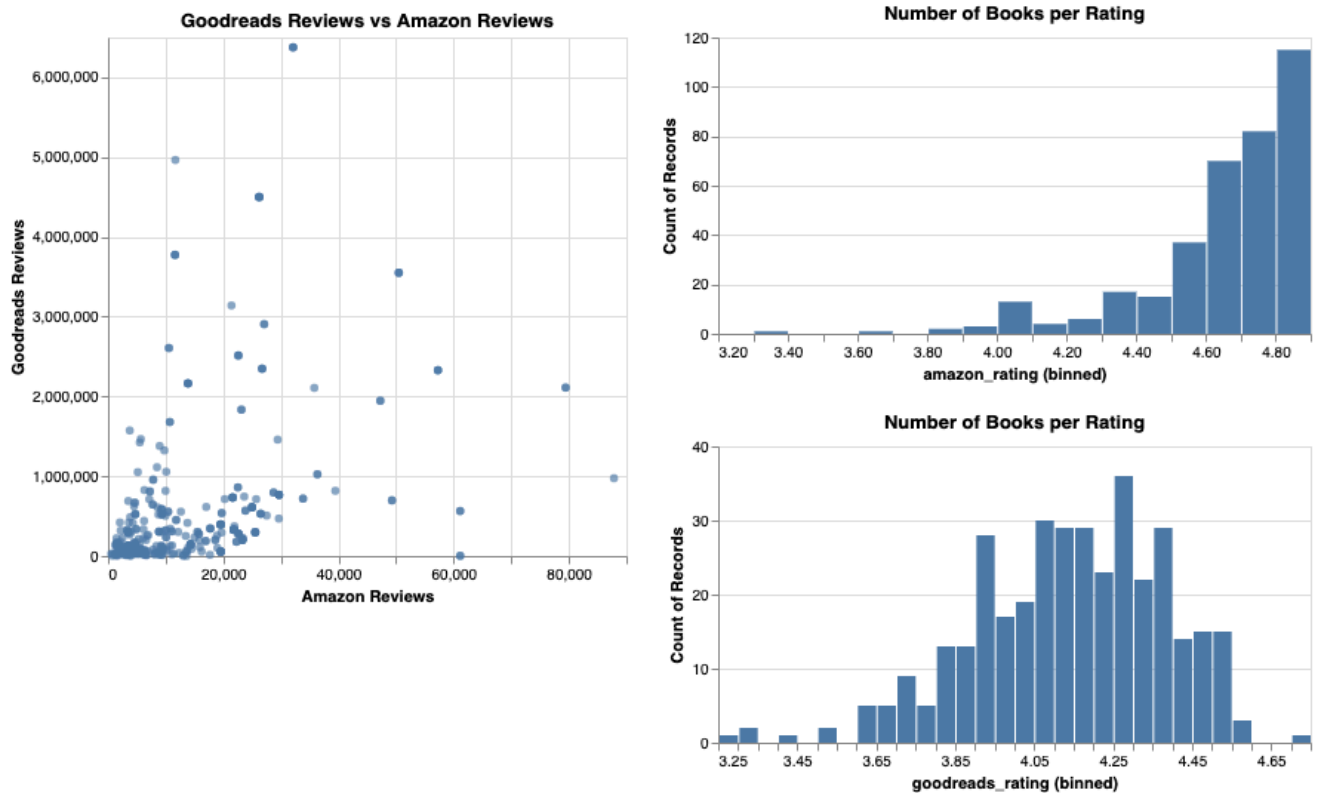
- Characteristics of Interaction and interactivity
  - Hover
    - Action: Select
    - Event type: hover

- - - Reaction type: highlight
    - Interaction coupling: unidirectional
    - Action Elements - Focus: direct, Presence: implicit, Granularity: atomic
    - Reaction Elements - Activation: Immediate, Spread: Self contained form, Flow: Discrete
  - Drop-down
    - Action: select
    - Event type: click
    - reaction type: change data
    - Interaction Coupling: unidirectional
    - Action Elements - Focus: indirect, Presence: Explicit, Granularity: Composite
    - Reaction Elements - Activation: Immediate, Spread: self-contained form, Flow: discrete
    - Affordances: upside down triangle indicate a drop-down menu with choices available to pick from
- Critique
  - The price differences are on unaligned axes, which have lower accuracy compared to an aligned axis. Comparing the price differences of books may be more difficult.
  - The line between amazon and goodreads can be perceived as connected and shows there is a relationship between the two points
  - Encoding the two different genres with color and shape may be redundant, but it doesn't use up channels for it to affect other attributes. It also sends a stronger message that they are two different sites.

View 6:

Task: What is the number of books per rating filtered by number of reviews on Amazon and Goodreads?

- Image of visualization



Include screenshots of previous iterations of the visualization

N/A. Visualization met expectations at first iteration.

- Task addressed: The three visualization answers to the task are the number of books per rating, filtered by the number of reviews on Amazon and Goodreads.
- Explanation:
  - Marks

    Scatter plot: circles - representing books.

    Bar chart (count of books per Amazon rating): bars - representing count of books per rating on Amazon.

    Bar chart (count of books per Goodreads rating): bars - representing count of books per rating on Goodreads.

○ Channels

Scatter plot: Goodreads reviews is encoded on y-axis and Amazon reviews is encoded on x-axis. Book titles is encoded as tooltip.

Bar chart (count of books per Amazon rating): Amazon rating is encoded on x-axis and count of books is encoded on y-axis.

Bar chart (count of books per Goodreads rating): Goodreads rating is encoded on x-axis and count of books is encoded on y-axis.

○ Characteristics of Channels that were exploited

Scatter plot: For x-axis, it runs from 0 to 90000 (reviews) and for y-axis, it runs from 0 to ~6500000 (reviews). With the use of common axis, it is easy for users to navigate/find books with needed number of Amazon reviews and Goodreads review.

Bar chart (count of books per Amazon rating):  For x-axis, it runs from ~3.2 to ~4.8 and for y-axis, it runs from 0 to 120 (count of books).

Bar chart (count of books per Goodreads rating):  For x-axis, it runs from ~3.25 to ~4.65 and for y-axis, it runs from 0 to 40 (count of books).
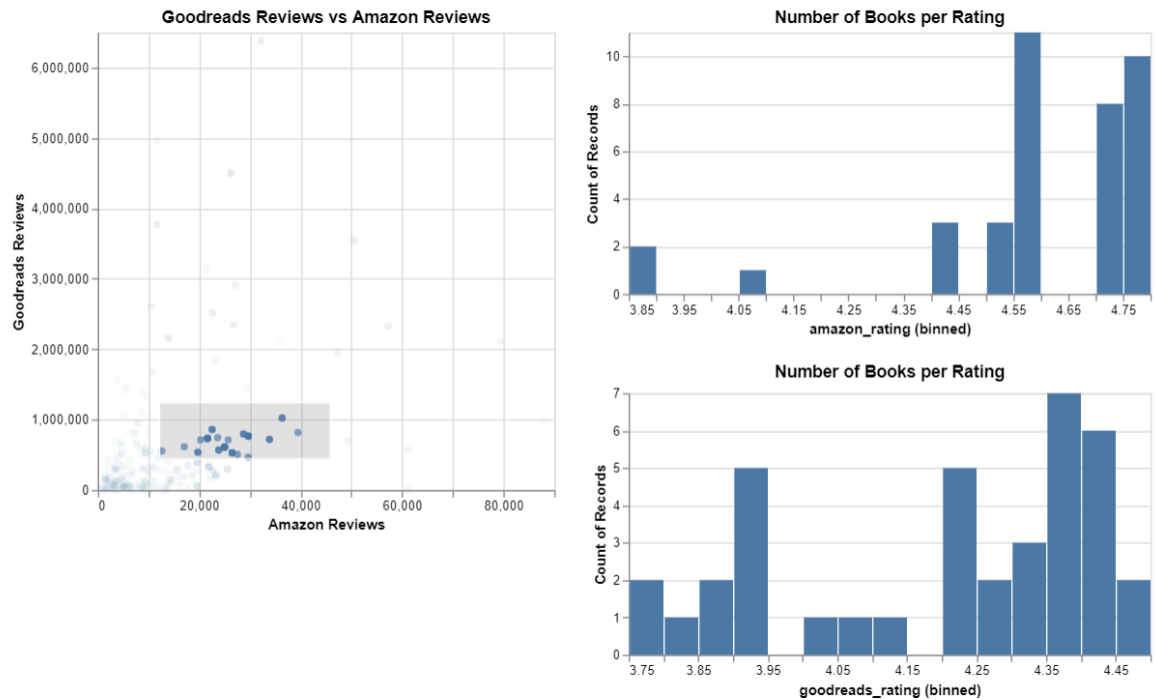
Since the bar chart uses common axis, it increases the accuracy of the visualization in terms of comparison of count of books within rating by comparing the height of bars to one another. .

○ Describe the interaction

Selection & filter (constrained on x-channel and y-channel of scatter plot, type: on shift-click, on mouse scroll, on click, semantics: null, group)

Users use interval selection on scatter plot to highlight the circles chosen (by changing the visual encoding) and filter the contents of bar charts. This allows user to see the count of books per Amazon rating and Goodreads rating for the selected data points.

Goodreads Reviews vs Amazon Reviews

Number of Books per Rating

Number of Books per Rating

○ Characteristics of Interaction and interactivity

Action: select & filter

Event type: on shift-click, on mouse scroll, on click

Reaction type: highlight, filter data

Views: juxtapose

Spread: propagated form - bidirectional. (affect data of bar chart only, scatter plot is unaffected)

Views data: share encoding, share data.

Interactivity: Focus: direct focus. Presence: implicit. Granularity: atomic

Activation: immediate

Flow: discrete flow.

○ Critique

Since bar charts use common axis, it increases the accuracy of the visualizations. Moreover, the scatter plot also uses common axis, so users can

navigate through x-axis and y-axis to find books with a specific number of Amazon reviews and Goodreads reviews.
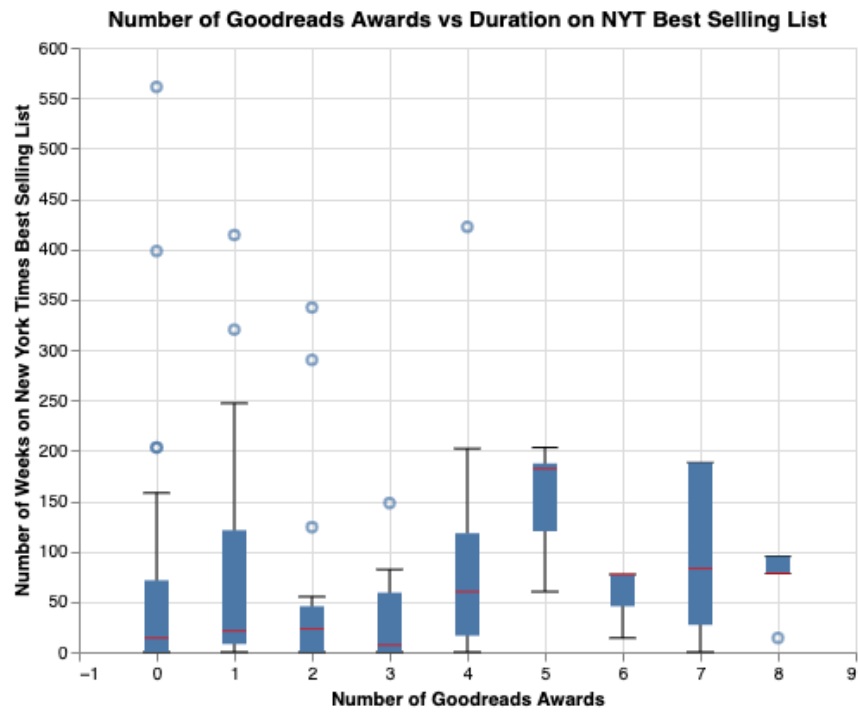
The use of interval selection on scatter plot provides user with freedom to filter the contents of bar charts based on number of Amazon reviews and Goodreads reviews. The charts together with the interaction creates juxtapose view that answers the task "What is the number of books per rating filtered by number of reviews on Amazon and Goodreads". The choice of charts are easy to understand for a wide range of audiences. Moreover, with the use of interactions in the view, users are provided with more freedom to filter data and encouraged to engage in the data analysis process.

On the other hand, this view could be improved by adding zooming on the scatter plot as most of the data points are concentrated in one area, which makes it harder to do a smaller range selection.
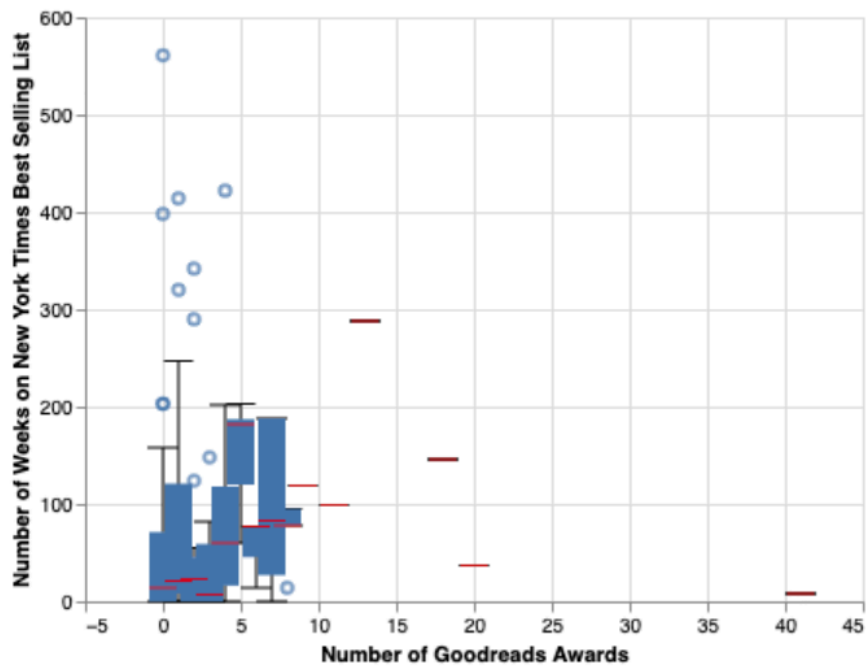
View 7:

Task: Do books with more awards tend to stay longer on the NYT best selling list?
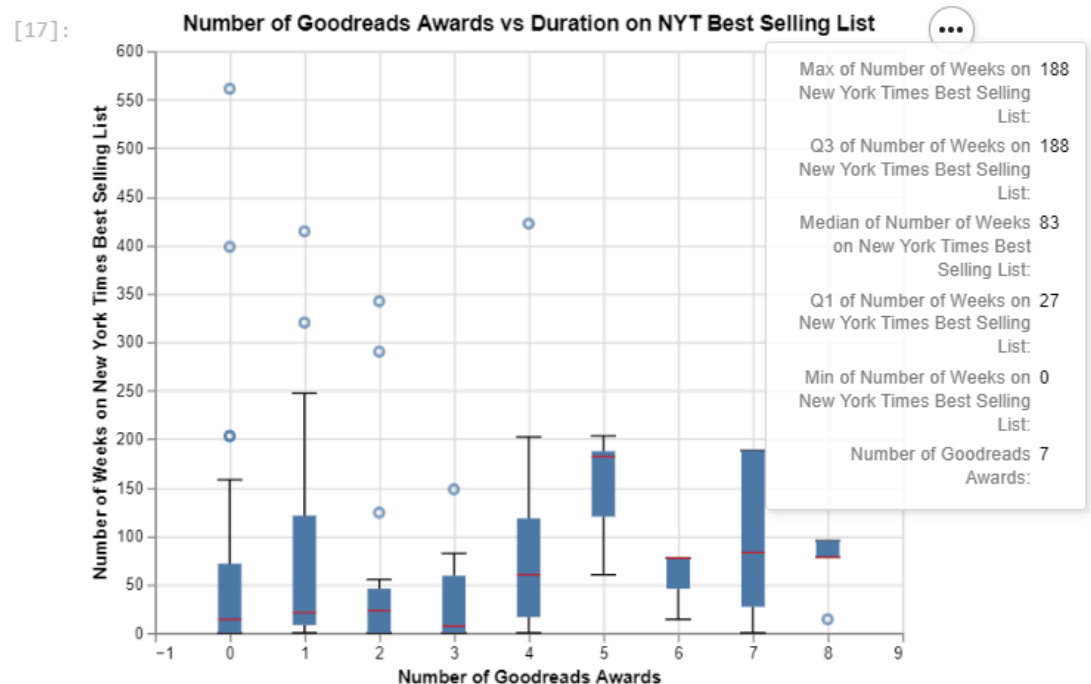- Image of visualization



- Include screenshots of previous iterations of the visualization

- Task addressed: Do books with more awards tend to stay longer on the NYT best-selling list?
- Explanation:
  - Marks:
    - box at each number of goodreads awards won
    - Points for outliers
  - Channels
    - Position on an unaligned axis to show the spread of the data
    - Horizontal position for the number of goodreads awards won
    - Tooltip encoded to express the 5 number summary of number of weeks on NYT best-selling list for each number of Goodreads awards.
  - Characteristics of Channels that were exploited
    - We want to compare
  - Describe the interaction
    - Tooltip states the 5 number summary of number of weeks on NYT best-selling list for each number of Goodreads awards.



  - Characteristics of Interaction and interactivity
    - Action: select
    - Event type: hover
    - Reaction type: highlight
    - Interaction coupling: unidirectional
    - Action Elements - Focus: direct, Presence: implicit, Granularity: atomic

- - - ■ Reaction Elements - Activation: Immediate, Spread: Self contained form, Flow: Discrete
  - ○ Critique
    - ■ Using boxplot gives a good summary of the data, but if we want to explore more complex questions we would probably need another view
    - ■ With the use of proximity, the levels of attribute: number of goodreads awards can be easily distinguished, which increases discriminability of channel x-axis.
    - ■ Position on an unaligned axis is not as accurate as position on an aligned axis. Comparing the spread between boxes may not be ideal in terms of accuracy.
    - ■ Popout with the red line to highlight the median which makes it easier to see

Novel Viz

- In addition to the visualizations that are programmed, you will be required to submit one high-fidelity sketch (i.e., show all the data) of a novel (must not be an existing viz idiom) visualization

Novel Viz



# Books Bestseller Hall of Fame

**Highest rating genre**

Self-improvement

★★★★ 
**700,000 reviews**

**MOST PROLIFIC PUBLISHER**

Penguin Random House

# of Bestseller books published

**28**

**Which book had the longest tenure on the Bestseller list?**

Tiny Changes, Remarkable Results
**Atomic Habits**
An Easy & Proven Way to Build Good Habits & Break Bad Ones
James Clear

Duration

**49 weeks**

**Author with most books in the Bestseller list**

Tien Nguyen

# OF BOOKS IN THE BESTSELLER LIST

**17**